Stanford Machine Learning w6 Advice.

## 1. Deciding What To do Next.
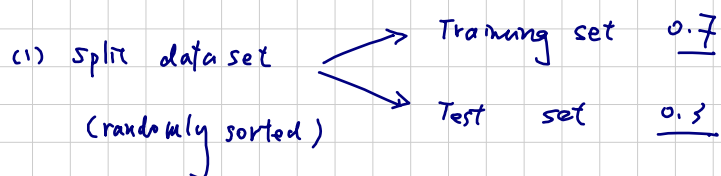
- Debugging → error ↑

  . larger training data.
  . smaller set of feature.
  . get additional features.
  . polynomial features.
  . ↑ $\lambda$   or   ↓ $\lambda$

  } how to pick one of the options?

- Machine Learning Diagnostic.

## 2. Evaluating a Hypothesis

(1) split data set

(randomly sorted)

→ Training set   $0.7$

→ Test   set   $0.3$    $M_{test}$ = # of test examples,

(2)   Procedure

(1) learn $\theta$ s.t. $\min_{\theta} J(\theta)$

(2) lin. reg.   $J_{test}(\theta) = \dfrac{1}{2M_{test}} \sum_{i=1}^{M_{test}} \left[ h_\theta(x_{test}^{(i)}) - y_{test}^{(i)} \right]^2$   (linear reg.)

log. reg   $J_{test}(\theta) = \dfrac{-1}{M_{test}} \sum_{i=1}^{M_{test}} y_{test}^{(i)} \cdot \log h_\theta(x_{test}^{(i)}) + (1-y_{test}^{(i)}) \log h_\theta(x_{test}^{(r)})$

misclassification error.

$err(h_\theta(x), y) = \begin{cases} 1 & (h_\theta(x) \geq 0.5 \;\&\; y=0) \;||\; (h_\theta(x) \leq 0.5 \;\&\; y=1) \\ 0 & otherwise \end{cases}$

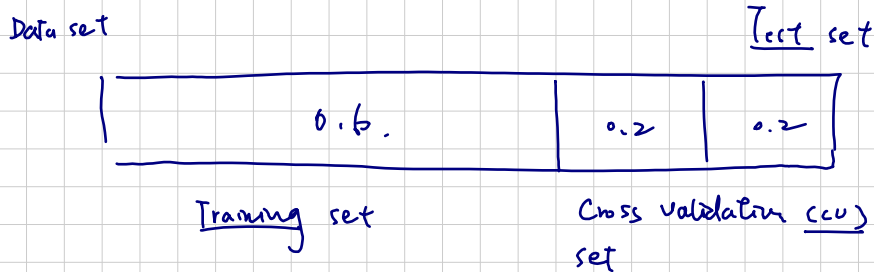Test error $= \dfrac{1}{M_{test}} \sum_{i=1}^{M_{test}} err(h_\theta(x_{test}^{(i)}), y_{test}^{(i)})$

# 3. Model Selection and Train / Validation / Test set

· Overfitting : cost < generalization error.

· Model selection ( $d = \deg_x(h)$ )

$$d = 1 \qquad h = \theta_0 + \theta_1 x \qquad \rightarrow \theta^{(1)} \rightarrow J_{test}(\theta^{(1)})$$
$$d = 2 \qquad h = \theta_0 + \theta_1 x + \theta_2 x^2 \qquad \rightarrow \theta^{(2)}$$
$$\vdots \qquad\qquad \vdots \qquad\qquad\qquad\qquad\qquad\qquad \vdots$$
$$d = 10 \qquad h = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10}. \rightarrow \theta^{(10)} \qquad J_{test}(\theta^{(10)})$$

· problem.  $J_{test}(\theta)$ is likely to be an optimistic generalization error.

(our extra parameter $d$ is fit to test)

Data set                                                    Test set

|  |  |  |
|---|---|---|
| 0.6. | 0.2 | 0.2 |

Training set        Cross validation (cv) set

Train / validation / test error.

Training error.   $J_{train}(\theta) = \dfrac{1}{2m} \sum\limits_{i=1}^{m} [h_\theta(x^{(i)}) - y^{(i)}]^2$

Cross validation error.   $J_{train}(\theta) = \dfrac{1}{2m_{cv}} \sum\limits_{i=1}^{m_{cv}} [h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)}]^2$

Test error   $J_{test}(\theta) = \dfrac{1}{2m_{test}} \sum\limits_{i=1}^{m_{test}} [h_\theta(x_{test}^{(i)}) - y_{test}^{(i)}]^2$

# Procedure

get $\theta^{(d)}$ → test on cv set. $J_{cv}(\theta^{(d)})$ $\forall d$.

→ find $d$ s.t. $J_{cv}(\theta^{(d)})$ is min.

→ estimate generalization error for test set $J_{test}(\theta^{(4)})$
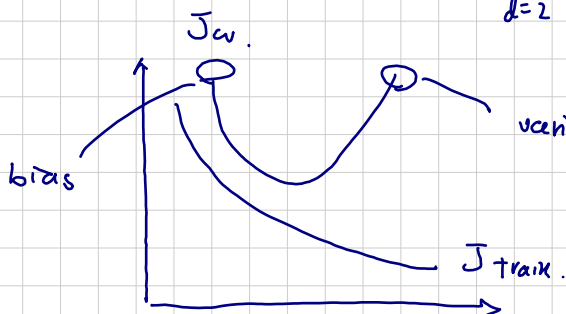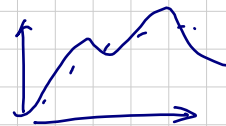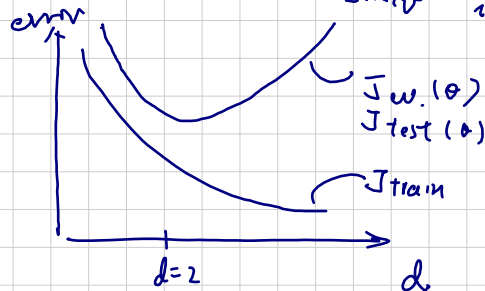
## 4. Diagnosing Bias v.s. Variance.



high bias
underfit

high variance
overfit

## Bias / Variance

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \left[ h_\theta(x^{(i)}) - y^{(i)} \right]^2$$

(or $J_{test}(\theta)$)

Cross validation error $\quad J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} \left[ h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)} \right]^2$



$J_{cv}(\theta)$
$J_{test}(\theta)$

$J_{train}$

$d=2$     $d$



$J_{cv}$

bias

variance (overfit. $d$ too large).

$J_{train}$.

⎡ Bias [underfit] → $J_{train}(\theta)$ high
⎢                  $J_{cv}(\theta) \approx J_{train}(\theta)$
⎢
⎣ Variance [overfit] → $J_{train}(\theta)$ low
                    $J_{cv}(\theta) >> J_{train}(\theta)$

# 5. Regularization and Bias / Variance.

Linear regression with regularization.

Model: $h_\theta(x) = \theta_0 + \theta_1 x_1 + \cdots + \theta_4 x^4$

$$J(\theta) = \underbrace{\frac{1}{2m} \sum_{i=1}^{m} [h_\theta(x^{(i)}) - y^{(i)}]^2}_{J_{train.} \ (w/o \ reg.)} + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2}_{reg.}$$

large $\lambda$.
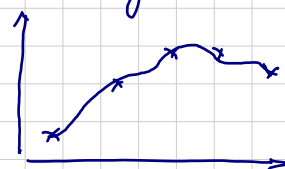
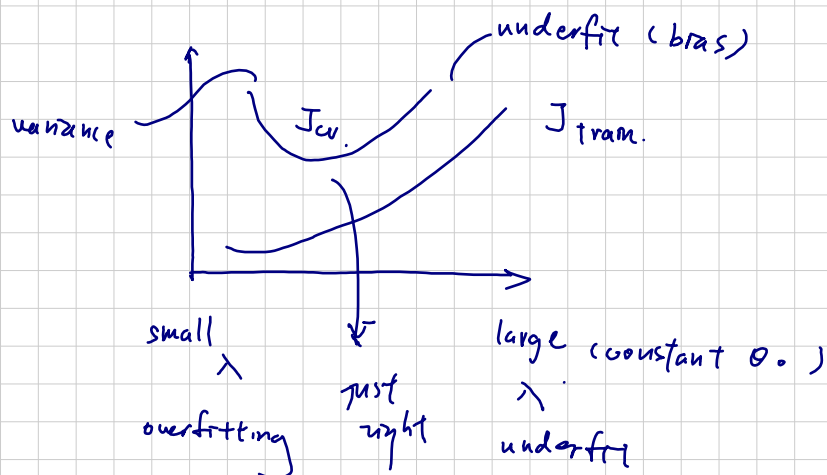$\lambda = 10000$, $\theta_{1,2,3,4} \doteq 0$.

$h_0 \simeq \theta_0$.

high bias (under fit)

$\lambda = 0$ (small $\lambda$)

high variance (overfit)

Training error. ($J_{train}$) vs. validation error ($J_{cv}$)

variance

$J_{cv}$.

underfit (bias)

$J_{train}$.

small $\lambda$

overfitting

just right

large (constant $\theta_0$.)

$\lambda$

underfit

# 6. Learning Curve.

the more data,
the better our
model.
generalizes

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} [h_\theta(x^{(i)}) - y^{(i)}]^2$$

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} [h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)}]^2$$

error. $\uparrow$ $J_{cv}(\theta)$

low error on
low m.

$m$ (training set size)

$\Rightarrow$ error $\sim m$. "easier to fit on

smaller training set"

for example $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$

$m=1$

$m=2$

$m=3$

## High Bias

error.
high $\updownarrow$

$J_{cv}(\theta)$ plateaus out.

$J_{train}(\theta) \approx J_{cv}$. (small data set)

$m$

If learning algorithm has high bias,
$\uparrow$ data set does not help.

<u>High Variance</u>

error.

$J_{cv}$.

Gap

$J_{train}$.

m

deg $h$. ↑
small $\lambda$.

If a learning algorithm is suffering from high variance, getting more
training data is likely to help.

7. Deciding What to Do Next Revisited.

· get more examples : fix high variance.
                                                        ⎫
· try smaller set of feature: fix high variance    ⎬ over fitting.
                                                        ⎭

· Add features . fix high bias.

· polynomial features : fix high bias

· try decr. $\lambda$ ⟶ fix high bias   ( incr. importance of $\theta$ ).

· try incr $\lambda$ ⟶ fix high variance (decr. importance of $\theta$ ).

<u>Neural Network and Overfitting</u>

"Small" Neural network
( few parameters ⟶ prone to
                          underfitting )

"Large" Neural Network
( more parameters ⟶ prone to overfitting )

computationally cheap.

how many
hidden layers?

Computationally        expensive.
⟶ use regularization (λ) to address
                          overfitting.

# Machine Learning System Design

1. Building a spam classifier

   · Supervised learning          $x$ = features of email.

                                  $y = 1$ (spam) or $0$ (not)

   · features $(x)$.    → indicative words.

   $$X = \begin{pmatrix} 1 \\ 0 \\ 1 \\ \vdots \\ \vdots \\ 1 \end{pmatrix} \begin{matrix} buy \\ discount \\ \vdots \end{matrix}$$

   presence of word.

   · Reduce error.
                    - acquire more data.
                    - develop sophisticated features based on email routing information
                      (from email header)
                    -              ''              for message body.
                    - develop sophisticated algorithm to detect misspellings
                      e.g.     w4tches

2. Error analysis.
              · Get something "quick and dirty" → plot learning curve
              (start w/ simple learning              → decide next steps
              algorithm) .

              · Error analysis ——→ manually look at examples and try to elicit systematic errors.

Ex.   $M_{cv} = 500$  examples.      Algorithm miss-classifies   100 emails .

                                                              ↳ manually examine.

(i)   what type of miss-classification

(ii)   what cues (features)  you think  would have helped   the

        algorithm  classify  them  correctly.


## Importance of  numerical evaluation

                    ↳ single  number  to   report  how  well  the algorithm is performing.

        - stemming  software .   ( discount / discounts / discounted )  ...  .

            (eg.  Porter Spammer )


## Skewed  Data

            - 1% error.   →  only   0.5% patients  have cancer.

            - skewed classes →    one class is  far more  ubiquitous

            - problem

                    - e.g.   99.2 %  accuracy    ( 0.8% error )

                                      ↓

                    99. 5%  accuracy    ( 0.5% error )

        for example,    function  y= predict cancer  (x)

                            y=0 ;

            return

                    → good accuracy  b/c of dataset

# 1. Precision / Recall.

Actual class

|  | 1 | 0 |
|---|---|---|
| prediction |  |  |
| 1 | (True positive) | false positive |
| 0 | false negative | True negative |

(1)  Precision  ( Of all patients where we predicted $y=1$, what fraction actually has cancer ? )

$$= \frac{\text{True positive}}{\# \text{ positive predictions}} = \frac{\text{true positives}}{(\text{true} + \text{false}) \text{ positives}}$$

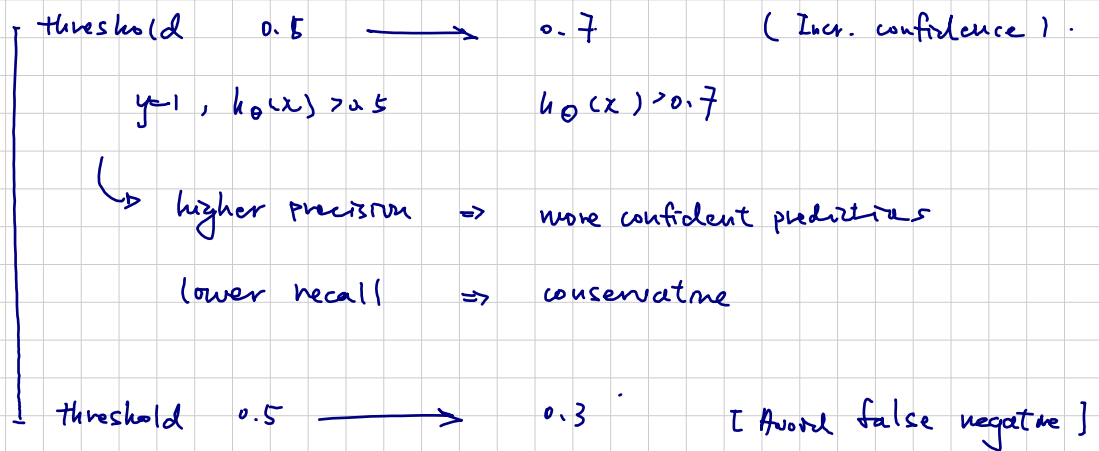(2)  Recall  ( of all patients that have cancer, what fraction did we predict so? )

$$= \frac{\text{true positive}}{\# \text{ actual positives}} = \frac{\text{true positive}}{\text{true positive} + \text{false negative}}$$

· Convention:  $y=1$  in presence of __rare class__

# 2. Trading off Precision and Recall.

$$\begin{cases} \text{precision} = \dfrac{\text{True pos.}}{\text{pred. pos}} \\[2em] \text{recall} = \dfrac{\text{true pos}}{\text{actual pos.}} \end{cases}$$
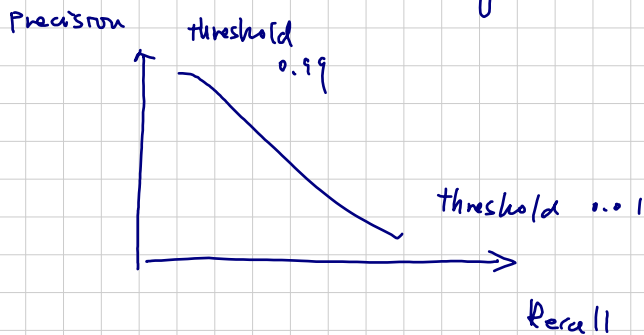
- logistic regression : $\quad 0 \le h_\theta(x) \le 1$

threshold    0.5   $\longrightarrow$   0.7      ( Incr. confidence ).

$y=1, \; h_\theta(x) > 0.5$      $h_\theta(x) > 0.7$

$\hookrightarrow$ higher precision $\Rightarrow$ more confident predictions

lower recall $\Rightarrow$ conservative

threshold    0.5   $\longrightarrow$   0.3      [ Avoid false negative ]

when in doubt, predict $y = 1$.

lower precision $\rightarrow$ less confident in our prediction.

higher recall. $\rightarrow$ less conservative



Precision

threshold 0.99

threshold 0.01

Recall

# How to select threshold?

. scalar that reports how well our algorithm performs.

$\hookrightarrow$ Now we have 2! ( precision / recall )

$\downarrow$ need a single number.

(1) Average. $\dfrac{P+R}{2}$ ?

- caveat :

(predict)

$y=1 >> y=0$     high recall.

$y=0 >> y=1$     high precision.

$\Big\}$ might skew our average

(2) F score.

$$F_1 = \dfrac{2\,PR}{P+R} \quad ; \quad P=0 \ \text{or} \ R=0 \ \longrightarrow \ F_1 = 0 .$$

$P=1$ and $R=1$ $\longrightarrow F_1 = 1$ ( perfect precision and recall ).

# Data for Machine Learning

- Algorithms:
    - perceptron.
    - Winnow.
    - memory-based
    - naïve bayes.

## Large Data rationale

- Assume feature $x \in \mathbb{R}^{n+1}$ has sufficient information to predict $y$ accurately.

- useful test: given the input $x$, can a human expert confidently predict $y$ ?

- Use learning algorithm with many parameters
    - e.g. logistic regression / linear regression with many features
      neural network with many hidden units.

        ↳ low bias algorithms.    $J_{train}(\theta)$ will be small

- Use large training set    (unlikely overfit).

                                    $J_{train} \approx J_{test}$

⇒    $J_{test}$ will be small